

# Mandoline: Model Evaluation under Distribution Shift

**Mayee Chen\***, Karan Goel\*, Nimit Sohoni\*, Fait Poms, Kayvon Fatahalian, Christopher Ré

*Stanford University*



# Motivation

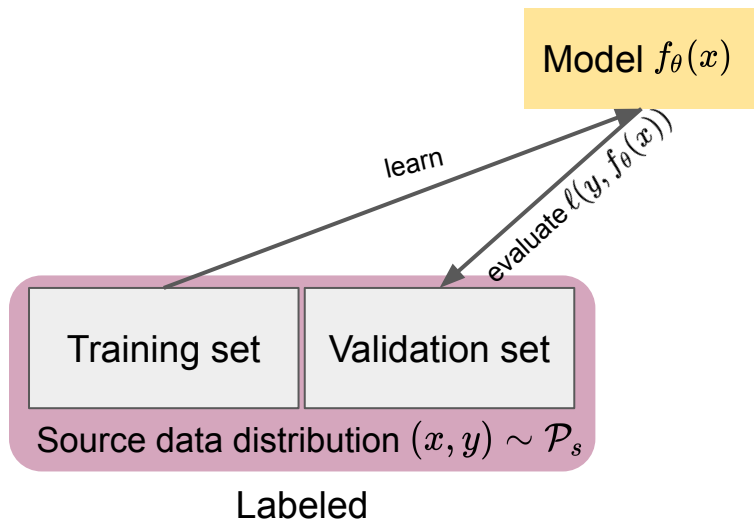
Q: How do we **evaluate** model performance during deployment?

- Model's deployment setting  $\neq$  training setting

# Motivation

Q: How do we **evaluate** model performance during deployment?

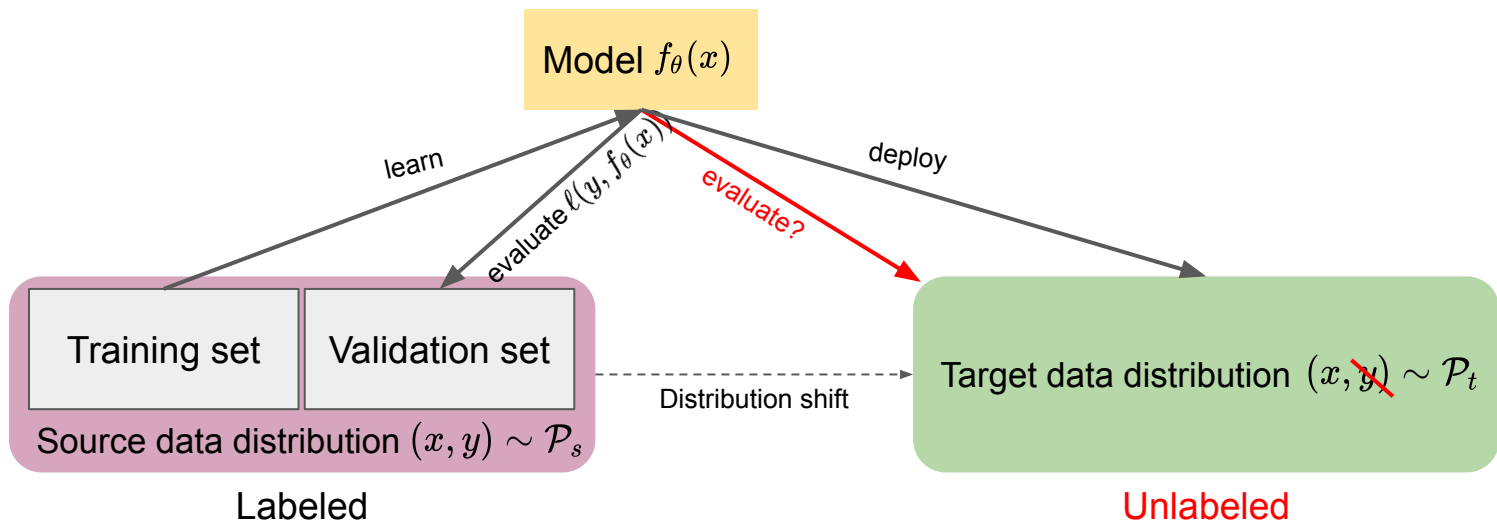
- Model's deployment setting  $\neq$  training setting



# Motivation

Q: How do we **evaluate** model performance during deployment?

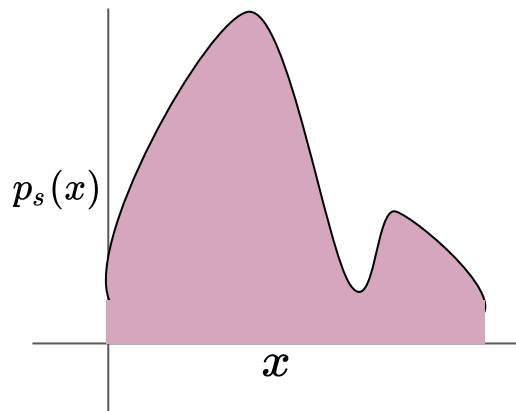
- Model's deployment setting  $\neq$  training setting



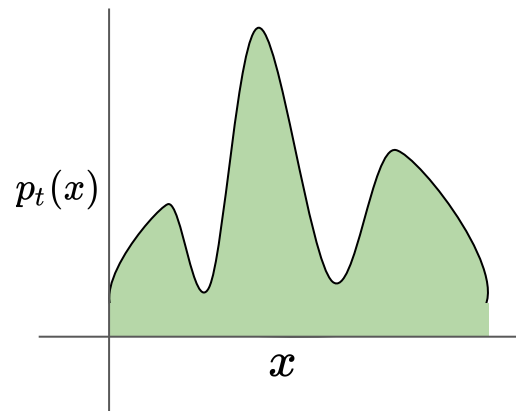
**Mandoline:** user-guided framework for evaluation under distribution shift

# Common approach: importance weighting

Source data distribution  $\mathcal{P}_s$



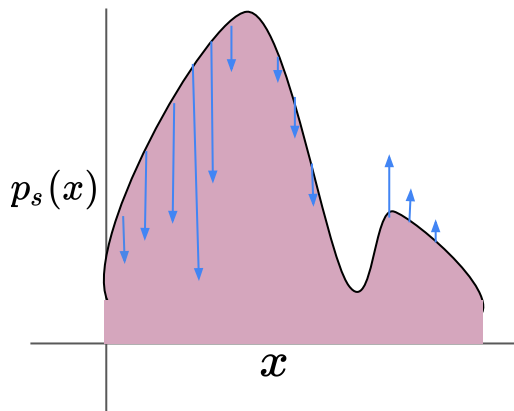
Target data distribution  $\mathcal{P}_t$



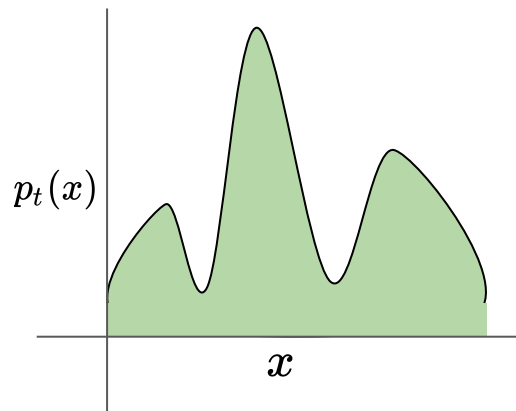
$$\mathbb{E}_t[\ell(y, f_\theta(x))] = \mathbb{E}_s\left[\frac{p_t(x)}{p_s(x)}\ell(y, f_\theta(x))\right] \approx \frac{1}{n} \sum_{i=1}^n \frac{p_t(x_i)}{p_s(x_i)}\ell(y_i, f_\theta(x_i))$$

# Common approach: importance weighting

Source data distribution  $\mathcal{P}_s$



Target data distribution  $\mathcal{P}_t$

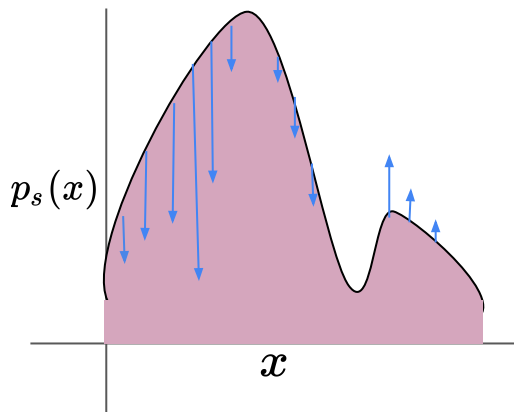


$$\mathbb{E}_t[\ell(y, f_\theta(x))] = \mathbb{E}_s\left[\frac{p_t(x)}{p_s(x)}\ell(y, f_\theta(x))\right] \approx \frac{1}{n} \sum_{i=1}^n \frac{p_t(x_i)}{p_s(x_i)} \ell(y_i, f_\theta(x_i))$$

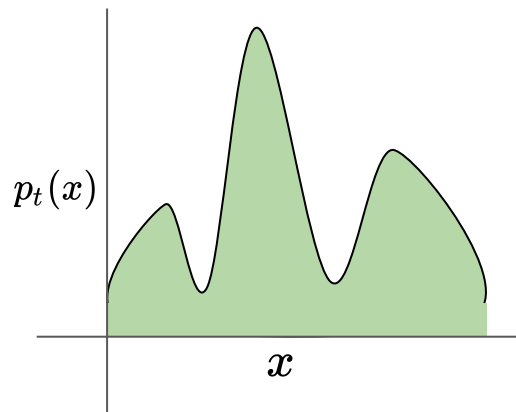
Density ratio to estimate

# Common approach: importance weighting

Source data distribution  $\mathcal{P}_s$



Target data distribution  $\mathcal{P}_t$



$$\mathbb{E}_t[\ell(y, f_\theta(x))] = \mathbb{E}_s\left[\frac{p_t(x)}{p_s(x)}\ell(y, f_\theta(x))\right] \approx \frac{1}{n} \sum_{i=1}^n \frac{p_t(x_i)}{p_s(x_i)} \ell(y_i, f_\theta(x_i))$$

Density ratio to estimate

Problems:

- Support shift - what if  $p_s(x) = 0, p_t(x) \neq 0$ ?
- High dimensional data  $x \in \mathbb{R}^d$ : harder to compute  $\frac{p_t(x)}{p_s(x)}$

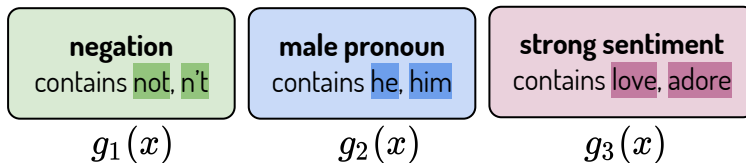
# Mandoline: Slice-based reweighting framework

***Slice***: user-defined grouping of data  $g(x) \in \{-1, 1\}$



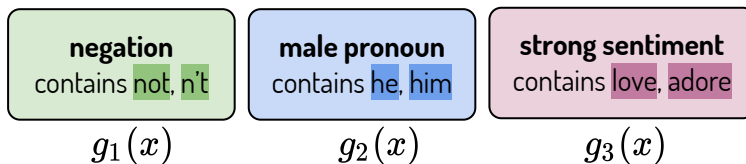
# Mandoline: Slice-based reweighting framework

**Slice:** user-defined grouping of data  $g(x) \in \{-1, 1\}$



# Mandoline: Slice-based reweighting framework

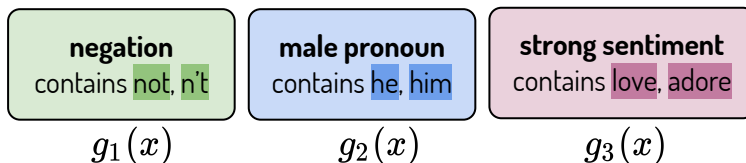
**Slice**: user-defined grouping of data  $g(x) \in \{-1, 1\}$



(Source) Labeled Validation Set	Slices			Model
I love eating ice-cream.	-1	-1	1	✓
He loved walking on the beach.	-1	1	1	✓
He didn't like drinking coffee.	1	1	-1	✗
⋮				
Source Accuracy: 91%				

# Mandoline: Slice-based reweighting framework

**Slice:** user-defined grouping of data  $g(x) \in \{-1, 1\}$



(Source) Labeled Validation Set	Slices			Model
I love eating ice-cream.	-1	-1	1	✓
He loved walking on the beach.	-1	1	1	✓
He didn't like drinking coffee.	1	1	-1	✗

⋮

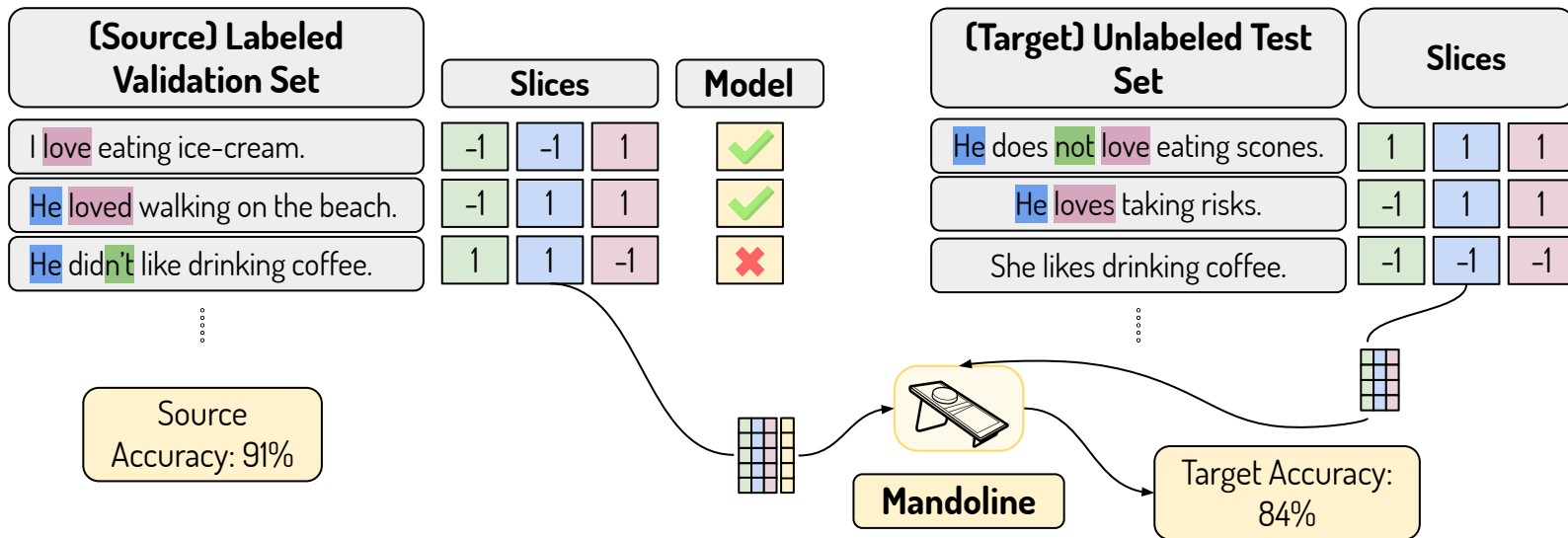
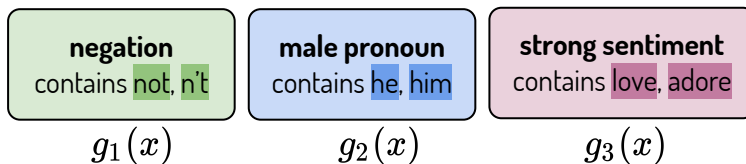
(Target) Unlabeled Test Set	Slices		
He does not love eating scones.	1	1	1
He loves taking risks.	-1	1	1
She likes drinking coffee.	-1	-1	-1

⋮

Source  
Accuracy: 91%

# Mandoline: Slice-based reweighting framework

**Slice**: user-defined grouping of data  $g(x) \in \{-1, 1\}$



# Results

**Prop 1:** if  $k$  slices  $g = \{g_1, \dots, g_k\}$  capture *all “relevant” distributional shift* between  $\mathcal{P}_s$  and  $\mathcal{P}_t$ , then reweighting with  $\frac{p_t(g(x))}{p_s(g(x))}$  recovers  $\mathbb{E}_t[\ell(y, f_\theta(x))]$ .

- If support shift occurs on irrelevant slices (i.e. slices independent of  $Y$ ), it can be corrected!
- Dimensionality: reduce from  $d \rightarrow k$

# Results

**Prop 1:** if  $k$  slices  $g = \{g_1, \dots, g_k\}$  capture *all “relevant” distributional shift* between  $\mathcal{P}_s$  and  $\mathcal{P}_t$ , then reweighting with  $\frac{p_t(g(x))}{p_s(g(x))}$  recovers  $\mathbb{E}_t[\ell(y, f_\theta(x))]$ .

- If support shift occurs on irrelevant slices (i.e. slices independent of  $Y$ ), it can be corrected!
- Dimensionality: reduce from  $d \rightarrow k$
- How to compute  $\frac{p_t(g(x))}{p_s(g(x))}$ ? Use any density ratio estimation method on  $g(x)$ 
  - **Kullback-Leibler Importance Estimation Procedure (KLIEP)**
    - Extend to correct for noisily-defined, incomplete slices

# Results

<b>Task</b>	<b>Task Labels</b>	<b>Distribution Shift</b>	<b>Slices</b>
CELEBA <i>image classification</i>	<i>male</i> <i>vs. female</i>	↑ blurry images	METADATA LABELS <i>blurry / not blurry</i>
SNLI→MNL <i>natural language</i> <i>inference</i>	<i>entailment, neutral</i> <i>or contradiction</i>	single-genre → multi-genre examples	PROGRAMMATIC <i>task model predictions,</i> <i>task model entropy</i>

# Results

Task	Task Labels	Distribution Shift	Slices
CELEBA <i>image classification</i>	<i>male</i> <i>vs. female</i>	↑ blurry images	METADATA LABELS <i>blurry / not blurry</i>
SNLI→MNL <i>natural language inference</i>	<i>entailment, neutral</i> <i>or contradiction</i>	single-genre → multi-genre examples	PROGRAMMATIC <i>task model predictions,</i> <i>task model entropy</i>

AVERAGE ESTIMATION ERROR (%)		
METHOD	CELEBA	
	RESNET18	RESNET50
SOURCE	1.96%	1.74%
CBIW	0.47%	0.53%
KMM	1.97%	1.76%
ULSIF	1.97%	1.76%
MANDOLINE	<b>0.16%</b>	<b>0.16%</b>

CelebA

Importance  
weighting on  $x$

On  $g(x)$

METHOD	STANDARD ACCURACY	
	AVG. ERROR	MAX. ERROR
SOURCE	6.2% ± 3.8%	15.6%
CBIW	5.5% ± 4.5%	17.9%
KMM	5.7% ± 3.6%	14.6%
ULSIF	6.4% ± 3.9%	16.0%
MANDOLINE	<b>3.6% ± 1.6%</b>	<b>5.9%</b>

SNLI → MNL



# Summary

Model evaluation under distribution shift:

- When user-specified slices capture relevant distribution shift, can reweight using them
- Can mitigate 1) support shift and 2) high dimensionality in standard importance reweighting
- Future steps: slice design - frameworks for how to construct good  $g$ ?



Thank you!

Contact: [mfchen@stanford.edu](mailto:mfchen@stanford.edu)